



Corrections to LRT on Large Dimensional Covariance Matrix by RMT

Zhidong Bai, Dandan Jiang, Jian-Feng Yao, Shurong Zheng

► To cite this version:

Zhidong Bai, Dandan Jiang, Jian-Feng Yao, Shurong Zheng. Corrections to LRT on Large Dimensional Covariance Matrix by RMT. *Annals of Statistics*, 2009, 37 (6B), pp.3822-3840. 10.1214/09-AOS694 . hal-00358579

HAL Id: hal-00358579

<https://hal.science/hal-00358579>

Submitted on 3 Feb 2009

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Corrections to LRT on Large Dimensional Covariance Matrix by RMT[†]

Zhidong Bai^{*} and Dandan Jiang and Jian-feng Yao and Shurong Zheng

Zhidong Bai, Dandan Jiang and Shurong Zheng
KLASMOE and School of Mathematics and Statistics
Northeast Normal University

5268 People's Road
130024 Changchun, China

and

Department of Statistics and Applied Probability
National University of Singapore
10, Kent Ridge Crescent
Singapore 119260

e-mail: stabaizd@nus.edu.sg, e-mail: stj@nus.edu.sg, e-mail: zhengsr@nenu.edu.cn

Jian-feng Yao
IRMAR and Université de Rennes 1
Campus de Beaulieu
35042 Rennes Cedex, France
e-mail: jian-feng.yao@univ-rennes1.fr

Abstract: In this paper, we give an explanation to the failure of two likelihood ratio procedures for testing about covariance matrices from Gaussian populations when the dimension is large compared to the sample size. Next, using recent central limit theorems for linear spectral statistics of sample covariance matrices and of random F-matrices, we propose necessary corrections for these LR tests to cope with high-dimensional effects. The asymptotic distributions of these corrected tests under the null are given. Simulations demonstrate that the corrected LR tests yield a realized size close to nominal level for both moderate p (around 20) and high dimension, while the traditional LR tests with χ^2 approximation fails.

Another contribution from the paper is that for testing the equality between two covariance matrices, the proposed correction applies equally for non-Gaussian populations yielding a valid pseudo-likelihood ratio test.

AMS 2000 subject classifications: Primary 62H15; secondary 62H10.

Keywords and phrases: High-dimensional data, Testing on covariance matrices, Marčenko-Pastrur distributions, Random F-matrices.

^{*}The research of this author was supported by CNSF grant 10571020 and NUS grant R-155-000-061-112

[†]This version contains full proofs of all results. A shorter version is to be published in a journal.

1. Introduction

The rapid development and wide application of computer techniques permits to collect and store a huge amount data, where the number of measured variables is usually large. Such high dimensional data occur in many modern scientific fields, such as micro-array data in biology, stock market analysis in finance and wireless communication networks. Traditional estimation or test tools are no more valid, or perform badly for such high-dimensional data, since they typically assume a large sample size n with respect to the number of variables p . A better approach in this high-dimensional data setting would be based on asymptotic theory which has both n and p approaching infinity. To illustrate this purpose, let us mention the case of Hotelling's T^2 -test. The failure of T^2 -test for high-dimensional data has been mentioned as early as by [Dempster \(1958\)](#). As a remedy, Dempster proposed a so-called non-exact test. However, the theoretical justification of Dempster's test arises much later in [Bai and Saranadasa \(1996\)](#) inspired by modern random matrix theory (RMT). These authors have found necessary correction for the T^2 -test to compensate effects due to high dimension.

In this paper, we consider two LR tests concerning covariance matrices. We first give a theoretical explanation for the fail of these tests in high-dimensional data context. Next, with the aid of random matrix theory, we provide necessary corrections to these LR tests to cope with the high dimensional effects.

First, we consider the problem of one-sample covariance hypothesis test. Suppose that \mathbf{x} follows a p -dimensional Gaussian distribution $N(\mu_p, \Sigma_p)$ and we want to test

$$H_0 : \Sigma_p = I_p , \quad (1.1)$$

where I_p denotes the p -dimensional identity matrix. Note that testing $\Sigma_p = A$ with an arbitrary covariance matrix A can always be reduced to the above null hypothesis by the transformation $A^{-\frac{1}{2}}\mathbf{x}$.

Let $(\mathbf{x}_1, \dots, \mathbf{x}_n)$ be a sample from \mathbf{x} , where we assume $p < n$. The sample covariance matrix is

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^p (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^*, \quad (1.2)$$

and set

$$L^* = \text{tr} \mathbf{S} - \log |\mathbf{S}| - p . \quad (1.3)$$

The likelihood ratio test statistic is

$$T_n = n \cdot L^* . \quad (1.4)$$

Keeping p fixed while letting $n \rightarrow \infty$, then the classical theory depicts that T_n converges to the $\chi_{\frac{1}{2}p(p+1)}^2$ distribution under H_0 .

However, as it will be shown, this classical approximation leads to a test size much higher than the nominal test level in the case of high-dimensional data, because T_n approaches infinity for large p . As seen from Table 1 in §3, for dimension and sample sizes $(p, n) = (50, 500)$, the realized size of the test is 22.5% instead of the nominal 5% level. The result is even worse for the case $(p, n) = (300, 500)$, with a 100% test size.

Based on a recent CLT for linear spectral statistics (LSS) of large-dimensional sample covariance matrices (Bai and Silverstein, 2004), we construct a corrected version of T_n in §3. As shown by the simulation results of §3.1, the corrected test performs much better in case of high dimensions. Moreover, it also performs correctly for moderate dimensions like $p = 10$ or 20 . For dimension and sample sizes (p, n) cited above, the sizes of the corrected test are 5.9% and 5.2%, respectively, both close to the 5% nominal level.

The second test problem we consider is about the equality between two high-dimensional covariance matrices. Let $\mathbf{x}_i = (x_{1i}, x_{2i}, \dots, x_{pi})^T, i = 1, \dots, n_1$ and $\mathbf{y}_j = (y_{1j}, y_{2j}, \dots, y_{pj})^T, j = 1, \dots, n_2$ be observations from two p -dimensional normal populations $N(\mu_k, \Sigma_k), k = 1, 2$, respectively. We wish to test the null hypothesis

$$H_0 : \Sigma_1 = \Sigma_2. \quad (1.5)$$

The related sample covariance matrices are

$$A = \frac{1}{n_1} \sum_{i=1}^{n_1} (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^*, \quad B = \frac{1}{n_2} \sum_{j=1}^{n_2} (\mathbf{y}_j - \bar{\mathbf{y}})(\mathbf{y}_j - \bar{\mathbf{y}})^*,$$

where $\bar{\mathbf{x}}, \bar{\mathbf{y}}$ are the respective sample means. Let

$$L_1 = \frac{|A|^{\frac{n_1}{2}} \cdot |B|^{\frac{n_2}{2}}}{|c_1 A + c_2 B|^{\frac{N}{2}}}, \quad (1.6)$$

where $N = n_1 + n_2$ and c_k denote $\frac{n_k}{N}, k = 1, 2$. The likelihood ratio test statistic is

$$T_N = -2 \log L_1,$$

and when $n_1, n_2 \rightarrow \infty$, we get

$$T_N = -2 \log L_1 \Rightarrow \chi_{\frac{1}{2}p(p+1)}^2 \quad (1.7)$$

under H_0 . Of course, in this limit scheme, the data dimension p is held fixed.

However, employing this χ^2 limit distribution for dimensions like 30 or 40, increases dramatically the size of the test. For instance, simulations in §4.1 show that, for dimension and sample sizes $(p, n_1, n_2) = (40, 800, 400)$, the test size equals 21.2% instead of the nominal 5% level. The result is worse for the case of $(p, n_1, n_2) = (80, 1600, 800)$, leading to a 49.5% test size. The reason for this fail of classical LR test is the following. Modern RMT indicates that when both dimension and sample size are large, the likelihood ratio statistic T_N drifts to infinity almost surely. Therefore,

the classical χ^2 approximation leads to many false rejections of H_0 in case of high-dimensional data.

Based on recent CLT for linear spectral statistics of F -matrices from RMT, we propose a correction to this LR test in §4. Although this corrected test is constructed under the asymptotic scheme $n_1 \wedge n_2 \rightarrow +\infty$, $y_{n_1} = p/n_1 \rightarrow y_1 \in (0, 1)$, $y_{n_2} = p/n_2 \rightarrow y_2 \in (0, 1)$, simulations demonstrate an overall correct behavior including small or moderate dimensions p . For example, for the above cited dimension and sample sizes (p, n_1, n_2) , the sizes of the corrected test equal 5.6% and 5.2%, respectively, both close to the nominal 5% level.

Related works include Ledoit and Wolf (2002), Srivastava (2005) and Schott (2007). These authors propose several procedures in the high-dimensional setting for testing that i) a covariance matrix is an identity matrix, proportional to an identity matrix (sphericity) and is a diagonal matrix or ii) several covariance matrices are equal. These procedures have the following common feature: their construction involves some well-chosen distance function between the null and the alternative hypotheses and rely on the first two spectral moments, namely the statistics $\text{tr}S_k$ and $\text{tr}S_k^2$ from sample covariance matrices S_k . Therefore, the procedures proposed by these authors are different from the likelihood-based procedures we consider here. Another important difference concerns the Gaussian assumption on the random variables used in all these references. Actually, for testing the equality between two covariance matrices, the correction proposed in this paper applies equally for non-Gaussian and high-dimensional data leading to a valid pseudo-likelihood test.

The rest of the paper is organized as following. Preliminary and useful RMT results are recalled in §2. In §3 and §4, we introduce our results for the two tests above. Proofs and technical derivations are postponed to the last section.

2. Useful results from the random matrix theory

We first recall several results from RMT, which will be useful for our corrections to tests. For any $p \times p$ square matrix M with real eigenvalues (λ_i^M) , F_n^M denotes the empirical spectral distribution (ESD) of M , that is,

$$F_n^M(x) = \frac{1}{p} \sum_{i=1}^p \mathbf{1}_{\lambda_i^M \leq x}, \quad x \in \mathbb{R}.$$

We will consider random matrix M whose ESD F_n^M converges (in a sense to be precised) to a limiting spectral distribution (LSD) F^M . To make statistical inference about a parameter $\theta = \int f(x) dF^M(x)$, it is natural to use the estimator

$$\hat{\theta} = \int f(x) dF_n^M(x) = \frac{1}{p} \sum_{i=1}^p f(\lambda_i^M),$$

which is a so-called linear spectral statistic (LSS) of the random matrix M .

2.1. CLT for LSS of a high-dimensional sample covariance matrix

Let $\{\xi_{ki} \in \mathbb{C}, i, k = 1, 2, \dots\}$ be a double array of *i.i.d.* complex variables with mean 0 and variance 1. Set $\xi_i = (\xi_{1i}, \xi_{2i}, \dots, \xi_{pi})^T$, the vectors (ξ_1, \dots, ξ_n) is considered as an *i.i.d* sample from some p -dimensional distribution with mean 0_p and covariance matrix I_p . Therefore the sample covariance matrix is

$$S_n = \frac{1}{n} \sum_{i=1}^n \xi_i \xi_i^*. \quad (2.1)$$

For $0 < \theta \leq 1$, let $a(\theta) = (1 - \sqrt{\theta})^2$ and $b(\theta) = (1 + \sqrt{\theta})^2$. The Marčenko-Pastur distribution of index θ , denoted as F^θ , is the distribution on $[a(\theta), b(\theta)]$ with the following density function

$$g_\theta(x) = \frac{1}{2\pi\theta x} \sqrt{[b(\theta) - x][x - a(\theta)]}, \quad a(\theta) \leq x \leq b(\theta).$$

Let

$$y_n = \frac{p}{n} \rightarrow y \in (0, 1)$$

and F^y, F^{y_n} be the Marčenko-Pastur law of index y and y_n , respectively. Let \mathcal{U} be an open set of the complex plane, including $[I_{(0,1)}(y)a(y), b(y)]$, and \mathcal{A} be the set of analytic functions $f : \mathcal{U} \mapsto \mathbb{C}$. We consider the empirical process $G_n := \{G_n(f)\}$ indexed by \mathcal{A} ,

$$G_n(f) = p \cdot \int_{-\infty}^{+\infty} f(x) [F_n - F^{y_n}](dx), \quad f \in \mathcal{A}, \quad (2.2)$$

where F_n is the ESD of S_n . The following theorem will play a fundamental role in next derivations, which is a specialization of a general theorem from [Bai and Silverstein \(2004\)](#) (Theorem 1.1).

Theorem 2.1. Assume that $f_1, \dots, f_k \in \mathcal{A}$, and $\{\xi_{ij}\}$ are *i.i.d.* random variables, such that $E\xi_{11} = 0$, $E|\xi_{11}|^2 = 1$, $E|\xi_{11}|^4 < \infty$. Moreover, $\frac{p}{n} \rightarrow y \in (0, 1)$ as $n, p \rightarrow \infty$.

Then:

(i) *Real Case.* Assume $\{\xi_{ij}\}$ are real and $E(\xi_{11}^4) = 3$. Then the random vector $(G_n(f_1), \dots, G_n(f_k))$ weakly converges to a k -dimensional Gaussian vector with mean vector,

$$m(f_j) = \frac{f_j(a(y)) + f_j(b(y))}{4} - \frac{1}{2\pi} \int_{a(y)}^{b(y)} \frac{f_j(x)}{\sqrt{4y - (x - 1 - y)^2}} dx, \quad j = 1, \dots, k, \quad (2.3)$$

and covariance function

$$v(f_j, f_\ell) = -\frac{1}{2\pi^2} \oint \oint \frac{f_j(z_1) f_\ell(z_2)}{(\underline{m}(z_1) - \underline{m}(z_2))^2} d\underline{m}(z_1) d\underline{m}(z_2), \quad j, \ell \in \{1, \dots, k\} \quad (2.4)$$

where $\underline{m}(z) \equiv m_{\underline{F}^y}(z)$ is the Stieltjes Transform of $\underline{F}^y \equiv (1 - y)I_{[0, \infty)} + yF^y$. The contours in (2.4) are non overlapping and both contain the support of F^y .

(ii) *Complex Case.* Assume $\{\xi_{ij}\}$ are complex and $E\xi_{11}^2 = 0$, $E(|\xi_{11}|^4) = 2$. Then the conclusion of (i) also holds, except the mean vector is zero and the covariance function is half of the function given in (2.4).

It is worth noticing that Theorem 1.1 in [Bai and Silverstein \(2004\)](#) covers more general sample covariance matrices of form $S'_n = T_n^{1/2} S_n T_n^{1/2}$ where (T_n) is a given sequence of positive-definite Hermitian matrices. In the “white” case $T_n \equiv I$ as considered here, in a recent preprint [Pastur and Lytova \(2008\)](#), the authors offer a new extension of the CLT where the constraints $E|\xi_{11}|^4 = 3$ or 2, as stated above, are removed.

2.2. CLT for LSS of high-dimensional F matrix

Let $\{\xi_{ki} \in \mathbb{C}, i, k = 1, 2, \dots\}$ and $\{\eta_{kj} \in \mathbb{C}, j, k = 1, 2, \dots\}$ are two independent double arrays of *i.i.d.* complex variables with mean 0 and variance 1. Write $\xi_i = (\xi_{1i}, \xi_{2i}, \dots, \xi_{pi})^T$ and $\eta_j = (\eta_{1j}, \eta_{2j}, \dots, \eta_{pj})^T$. Also, for any positive integers n_1, n_2 , the vectors $(\xi_1, \dots, \xi_{n_1})$ and $(\eta_1, \dots, \eta_{n_2})$ can be thought as independent samples of size n_1 and n_2 , respectively, from some p -dimensional distributions. Let S_1 and S_2 be the associated sample covariance matrices, *i.e.*

$$S_1 = \frac{1}{n_1} \sum_{i=1}^{n_1} \xi_i \xi_i^* \quad \text{and} \quad S_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} \eta_j \eta_j^*$$

Then, the following so-called F-matrix generalizes the classical Fisher-statistics for the present p -dimensional case,

$$V_n = S_1 S_2^{-1} \tag{2.5}$$

where $n_2 > p$. Here we use the notation $n = (n_1, n_2)$.

Let

$$y_{n_1} = \frac{p}{n_1} \rightarrow y_1 \in (0, 1), \quad y_{n_2} = \frac{p}{n_2} \rightarrow y_2 \in (0, 1). \tag{2.6}$$

Under suitable moment conditions, the ESD $F_n^{V_n}$ of V_n has a LSD F_{y_1, y_2} , which has a density [See P72 of [Bai and Silverstein \(2006\)](#)], given by

$$\ell(x) = \begin{cases} \frac{(1-y_2)\sqrt{(b-x)(x-a)}}{2\pi x(y_1+y_2x)}, & a \leq x \leq b, \\ 0, & \text{otherwise.} \end{cases} \tag{2.7}$$

where $a = (1-y_2)^{-2}(1-\sqrt{y_1+y_2-y_1y_2})^2$ and $b = (1-y_2)^{-2}(1+\sqrt{y_1+y_2-y_1y_2})^2$.

Similar to previously, let $\tilde{\mathcal{U}}$ be an open set of the complex plane, including the interval

$$\left[I_{(0,1)}(y_1) \frac{(1-\sqrt{y_1})^2}{(1+\sqrt{y_2})^2}, \frac{(1+\sqrt{y_1})^2}{(1-\sqrt{y_2})^2} \right],$$

and $\tilde{\mathcal{A}}$ be the set of analytic functions $f : \tilde{\mathcal{U}} \mapsto \mathbb{C}$. Define the empirical process $\tilde{G}_n := \{\tilde{G}_n(f)\}$ indexed by $\tilde{\mathcal{A}}$

$$\tilde{G}_n(f) = p \cdot \int_{-\infty}^{+\infty} f(x) [F_n^{V_n} - F_{y_{n_1}, y_{n_2}}](dx), \quad f \in \tilde{\mathcal{A}}. \tag{2.8}$$

Here $F_{y_{n_1}, y_{n_2}}$ is the limiting distribution in (2.7) but with y_{n_k} instead of $y_k, k = 1, 2$.

Recently, [Zheng \(2008\)](#) establishes a general CLT for LSS of large-dimensional F matrix. The following theorem is a simplified one quoted from it, which will play an important role.

Theorem 2.2. *Let $f_1, \dots, f_k \in \tilde{\mathcal{A}}$, and assume:*

For each p , (ξ_{ij_1}) and (η_{ij_2}) variables are i.i.d., $1 \leq i \leq p$, $1 \leq j_1 \leq n_1$, $1 \leq j_2 \leq n_2$. $E\xi_{11} = E\eta_{11} = 0$, $E|\xi_{11}|^4 = E|\eta_{11}|^4 < \infty$, $y_{n_1} = \frac{p}{n_1} \rightarrow y_1 \in (0, 1)$, $y_{n_2} = \frac{p}{n_2} \rightarrow y_2 \in (0, 1)$.

Then

(i) Real Case. Assume (ξ_{ij}) and (η_{ij}) are real, $E|\xi_{11}|^2 = E|\eta_{11}|^2 = 1$, then the random vector $(\widetilde{G}_n(f_1), \dots, \widetilde{G}_n(f_k))$ weakly converges to a k -dimensional Gaussian vector with the mean vector

$$m(f_j) = \lim_{r \rightarrow 1+} [(\textcolor{red}{2.9}) + (\textcolor{red}{2.10}) + (\textcolor{red}{2.11})]$$

$$\frac{1}{4\pi i} \oint_{|\zeta|=1} f_j(z(\zeta)) \left[\frac{1}{\zeta - \frac{1}{r}} + \frac{1}{\zeta + \frac{1}{r}} - \frac{2}{\zeta + \frac{y_2}{hr}} \right] d\zeta \quad (2.9)$$

$$+ \frac{\beta \cdot y_1(1-y_2)^2}{2\pi i \cdot h^2} \oint_{|\zeta|=1} f_j(z(\zeta)) \frac{1}{(\zeta + \frac{y_2}{hr})^3} d\zeta \quad (2.10)$$

$$+ \frac{\beta \cdot y_2(1-y_2)}{2\pi i \cdot h} \oint_{|\zeta|=1} f_j(z(\zeta)) \frac{\zeta + \frac{1}{hr}}{(\zeta + \frac{y_2}{hr})^3} d\zeta, \quad j = 1, \dots, k, \quad (2.11)$$

where $z(\zeta) = (1-y_2)^{-2} [1 + h^2 + 2h\mathcal{R}(\zeta)]$, $h = \sqrt{y_1 + y_2 - y_1 y_2}$, $\beta = E|\xi_{11}|^4 - 3$, and the covariance function as $1 < r_1 < r_2 \downarrow 1$

$$v(f_j, f_\ell) = \lim_{1 < r_1 < r_2 \rightarrow 1+} [(\textcolor{red}{2.12}) + (\textcolor{red}{2.13})]$$

$$- \frac{1}{2\pi^2} \oint_{|\zeta_2|=1} \oint_{|\zeta_1|=1} \frac{f_j(z(r_1\zeta_1)) f_\ell(z(r_2\zeta_2)) r_1 r_2}{(r_2\zeta_2 - r_1\zeta_1)^2} d\zeta_1 d\zeta_2, \quad (2.12)$$

$$- \frac{\beta \cdot (y_1 + y_2)(1-y_2)^2}{4\pi^2 h^2} \oint_{|\zeta_1|=1} \frac{f_j(z(\zeta_1))}{(\zeta_1 + \frac{y_2}{hr_1})^2} d\zeta_1 \oint_{|\zeta_2|=1} \frac{f_\ell(z(\zeta_2))}{(\zeta_2 + \frac{y_2}{hr_2})^2} d\zeta_2 \quad (2.13)$$

$j, \ell \in \{1, \dots, k\}$.

(ii) Complex Case. Assume (ξ_{ij}) and (η_{ij}) are complex, $E(\xi_{11}^2) = E(\eta_{11}^2) = 0$, then the conclusion of (i) also holds, except the means are $\lim_{r \rightarrow 1+} [(\textcolor{red}{2.10}) + (\textcolor{red}{2.11})]$ and the covariance function is

$$\lim_{1 < r_1 < r_2 \rightarrow 1+} \left[\frac{1}{2} \cdot (\textcolor{red}{2.12}) + (\textcolor{red}{2.13}) \right], \text{ where } \beta = E|\xi_{11}|^4 - 2.$$

We should point out that Zheng's CLT for F -matrices covers more general situations than those cited in Theorem 2.2. In particular, the fourth-moments $E|\xi_{11}|^4$ and $E|\eta_{11}|^4$ can be different.

The following lemma will be used in §4 for an application of Theorem 2.2 to obtain the formula (4.5) and (4.6).

Lemma 2.1. *For the function $f(x) = \log(a + bx)$, $x \in \mathbb{R}$, $a, b > 0$, let (c, d) be the unique solution to the equations*

$$\begin{cases} c^2 + d^2 = a(1-y_2)^2 + b(1+h^2), \\ cd = bh, \\ 0 < d < c. \end{cases}$$

Analogously, let γ, η be the constants similar to (c, d) but for the function $g(x) = \log(\alpha + \beta x)$, $\alpha > 0$, $\beta > 0$. Then, the mean and covariance functions in (2.9) and (2.12) equal to

$$\begin{aligned} m(f) &= \frac{1}{2} \log \frac{(c^2 - d^2)h^2}{(ch - y_2 d)^2}, \\ v(f, g) &= 2bhd^{-1}c^{-1} \log \frac{c\gamma}{c\gamma - d\eta}. \end{aligned}$$

3. Testing the hypothesis that a high-dimensional covariance matrix is equal to a given matrix

To test the hypothesis $H_0 : \Sigma_p = I_p$, let be the sample covariance matrix \mathbf{S} and likelihood ratio statistic T_n as defined in (1.2) and (1.4), respectively. For $\xi_i = \mathbf{x}_i - \mu_p$, the array $\{\xi_i\}_{i=1, \dots, n}$ contains p -dimensional standard normal variables under H_0 . Let

$$\mathbf{S}_n = \frac{1}{n} \sum_{i=1}^n \xi_i \xi_i^*.$$

and

$$\widetilde{L}^* = \text{tr} \mathbf{S}_n - \log |\mathbf{S}_n| - p.$$

Theorem 3.1. *Assuming that the conditions of Theorem 2.1 hold, L^* is defined as (1.3) and $g(x) = x - \log x - 1$. Then, under H_0 and when $n \rightarrow \infty$*

$$\widetilde{T}_n = v(g)^{-\frac{1}{2}} [L^* - p \cdot F^{y_n}(g) - m(g)] \Rightarrow N(0, 1), \quad (3.1)$$

where F^{y_n} is the Marčenko-Pastur law of index y_n .

Proof. Because the difference between \mathbf{S} and \mathbf{S}_n is a rank-1 matrix, \mathbf{S} and \mathbf{S}_n have the same LSD. So, L^* and \widetilde{L}^* have the same asymptotic distribution. We also have

$$\begin{aligned} \widetilde{L}^* &= \text{tr} \mathbf{S}_n - \log |\mathbf{S}_n| - p \\ &= \sum_{i=1}^p (\lambda_i^{\mathbf{S}_n} - \log \lambda_i^{\mathbf{S}_n} - 1) = p \cdot \int (x - \log x - 1) dF_n(x) \\ &= p \cdot \int g(x) d(F_n(x) - F^{y_n}(x)) + p \cdot F^{y_n}(g), \end{aligned}$$

so that

$$G_n(g) = \widetilde{L}^* - p \cdot F^{y_n}(g). \quad (3.2)$$

By Theorem 2.1, $G_n(g)$ weakly converges to a Gaussian vector with the mean

$$m(g) = -\frac{\log(1-y)}{2} \quad (3.3)$$

and variance

$$v(g) = -2 \log(1-y) - 2y. \quad (3.4)$$

for the real case, which are calculated in §5. For the complex case, the mean $m(g)$ is zero and the variance is half of $v(g)$. Then, by (3.2) we arrive at

$$\widetilde{L}^* - p \cdot F^{y_n}(g) \Rightarrow N(m(g), v(g)), \quad (3.5)$$

where

$$F^{y_n}(g) = 1 - \frac{y_n - 1}{y_n} \log(1 - y_n) \quad (3.6)$$

can be calculated by the density of LSD of sample covariance matrix in §5. Because \widetilde{L}^* and L^* have the same asymptotic distribution and (3.5), finally we get

$$\widetilde{T}_n = v(g)^{-\frac{1}{2}} [L^* - p \cdot F^{y_n}(g) - m(g)] \Rightarrow N(0, 1).$$

□

3.1. Simulation study I

For different values of (p, n) , we compute the realized sizes of traditional likelihood ratio test (LRT) and the corrected likelihood ratio test (CLRT) proposed previously. The nominal test level is set to be $\alpha = 0.05$, and for each (p, n) , we run 10,000 independent replications with real Gaussian variables. Results are given in Table 1 and Figure 1 below.

(p, n)	CLRT			LRT	
	Size	Difference with 5%	Power	Size	Power
(5, 500)	0.0803	0.0303	0.6013	0.0521	0.5233
(10, 500)	0.0690	0.0190	0.9517	0.0555	0.9417
(50, 500)	0.0594	0.0094	1	0.2252	1
(100, 500)	0.0537	0.0037	1	0.9757	1
(300, 500)	0.0515	0.0015	1	1	1

TABLE 1

Sizes and powers of the traditional LRT and the corrected LRT, based on 10,000 independent applications with real Gaussian variables. Powers are estimated under the alternative $\Sigma_p = \text{diag}(1, 0.05, 0.05, 0.05, \dots)$.

As seen from Table 1, the traditional LRT always rejects H_0 when p is large, like $p = 100$ or 300 , while the sizes produced by the corrected LRT perfectly matches the nominal level. For moderate dimensions like $p = 50$, the corrected LRT still performs correctly while the traditional LRT has a size much higher than 5%.

4. Testing the equality of two high-dimensional covariance matrices

Let (\mathbf{x}_i) , $i = 1, \dots, n_1$ and (\mathbf{y}_j) , $j = 1, \dots, n_2$ be observations from two normal populations $N(\mu_k, \Sigma_k)$, $k = 1, 2$, respectively. We examine the test defined in (1.5) and (1.6). The aim is to

find a good scaling of the LR statistic T_N , such that the scaled statistic weakly converges to some limiting distribution. Let

$$\xi_i = \Sigma^{-\frac{1}{2}}(\mathbf{x}_i - \mu_1), \quad \eta_i = \Sigma^{-\frac{1}{2}}(\mathbf{y}_i - \mu_2)$$

where $\Sigma = \Sigma_1 = \Sigma_2$ denotes the common covariance matrix under H_0 . Note that in a strict sense, the vectors $(\mathbf{x}_i), (\mathbf{y}_i)$ and the matrices $\Sigma, \Sigma_1, \Sigma_2$ depend on p . However we do not signify this dependence in notations for ease of statements. Due to Gaussian assumption, the arrays $(\xi_i)_{i=1, \dots, n_1}$ and $(\eta_j)_{j=1, \dots, n_2}$ contain i.i.d. $N(0, 1)$ variables, for which we can apply Theorem 2.2.

Let

$$\begin{aligned} S_1 &= \frac{1}{n_1} \sum_{i=1}^{n_1} \xi_i \xi_i^* = \Sigma^{-\frac{1}{2}} C \Sigma^{-\frac{1}{2}} \\ S_2 &= \frac{1}{n_2} \sum_{j=1}^{n_2} \eta_j \eta_j^* = \Sigma^{-\frac{1}{2}} D \Sigma^{-\frac{1}{2}}, \end{aligned}$$

where

$$\begin{aligned} C &= \frac{1}{n_1} \sum_{i=1}^{n_1} (\mathbf{x}_i - \mu_1)(\mathbf{x}_i - \mu_1)^*, \\ D &= \frac{1}{n_2} \sum_{j=1}^{n_2} (\mathbf{y}_j - \mu_2)(\mathbf{y}_j - \mu_2)^*. \end{aligned}$$

Note that

$$V_n = S_1 S_2^{-1}$$

forms a random F-matrix and we have

$$\widetilde{L}_1 = \frac{|S_1|^{\frac{n_1}{2}} \cdot |S_2|^{\frac{n_2}{2}}}{|c_1 S_1 + c_2 S_2|^{\frac{N}{2}}} = \frac{|C|^{\frac{n_1}{2}} \cdot |D|^{\frac{n_2}{2}}}{|c_1 C + c_2 D|^{\frac{N}{2}}}. \quad (4.1)$$

Theorem 4.1. *Assuming that the conditions of Theorem 2.2 hold under H_0 , L_1 as defined in (1.6) and*

$$f(x) = \log(y_{n_1} + y_{n_2}x) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log x - \log(y_{n_1} + y_{n_2}).$$

Then, under H_0 and as $n_1 \wedge n_2 \rightarrow \infty$,

$$\widetilde{T}_N = v(f)^{-\frac{1}{2}} \left[-\frac{2 \log L_1}{N} - p \cdot F_{y_{n_1}, y_{n_2}}(f) - m(f) \right] \Rightarrow N(0, 1). \quad (4.2)$$

Proof. As $A - C$ and $B - D$ are rank-1 random matrices, AB^{-1} and CD^{-1} have the same LSD.

Also by (4.1), \widetilde{L}_1 and L_1 have the same asymptotic distribution. Because

$$\begin{aligned} -\frac{2}{N} \log \widetilde{L}_1 &= -\frac{2}{N} \log \left(\frac{|S_1|^{\frac{n_1}{2}} \cdot |S_2|^{\frac{n_2}{2}}}{|c_1 S_1 + c_2 S_2|^{\frac{N}{2}}} \right) \\ &= \log |c_1 V_n^{-1} + c_2| - c_1 \cdot \log |V_n^{-1}| \\ &= \sum_{i=1}^p \log(c_1 \lambda_i^{V_n} + c_2) - c_1 \cdot \log(\lambda_i^{V_n}) \end{aligned}$$

$$= p \cdot \int [\log(c_1 x + c_2) - c_1 \cdot \log(x)] dF_n^{V_n}(x).$$

Define $f(x) = \log(c_1 x + c_2) - c_1 \cdot \log(x)$, by $c_1 = \frac{n_1}{N} = \frac{y_{n_2}}{y_{n_1} + y_{n_2}}$ and $c_2 = \frac{n_2}{N} = \frac{y_{n_1}}{y_{n_1} + y_{n_2}}$, also it can be written as

$$f(x) = \log(y_{n_1} + y_{n_2} x) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log x - \log(y_{n_1} + y_{n_2}). \quad (4.3)$$

From

$$\begin{aligned} -\frac{2 \log \widetilde{L}_1}{N} &= p \cdot \int f(x) dF_n^{V_n}(x) \\ &= p \cdot \int f(x) d(F_n^{V_n}(x) - F_{y_{n_1}, y_{n_2}}(x)) + p \cdot F_{y_{n_1}, y_{n_2}}(f), \end{aligned}$$

we get

$$\widetilde{G}_n(f) = -\frac{2 \log \widetilde{L}_1}{N} - p \cdot F_{y_{n_1}, y_{n_2}}(f). \quad (4.4)$$

By Theorem 2.2, $\widetilde{G}_n(f)$ weakly converges to a Gaussian vector with mean

$$m(f) = \frac{1}{2} \left[\log \left(\frac{y_1 + y_2 - y_1 y_2}{y_1 + y_2} \right) - \frac{y_1}{y_1 + y_2} \log(1 - y_2) - \frac{y_2}{y_1 + y_2} \log(1 - y_1) \right] \quad (4.5)$$

and variance

$$v(f) = -\frac{2y_2^2}{(y_1 + y_2)^2} \log(1 - y_1) - \frac{2y_1^2}{(y_1 + y_2)^2} \log(1 - y_2) - 2 \log \frac{y_1 + y_2}{y_1 + y_2 - y_1 y_2} \quad (4.6)$$

for the real case, which are calculated by Lemma 2.1 in §5. For the complex case, the mean $m(f)$ is zero and the variance is half of $v(f)$. In other words,

$$-\frac{2 \log \widetilde{L}_1}{N} - p \cdot F_{y_{n_1}, y_{n_2}}(f) \Rightarrow N(m(f), v(f)), \quad (4.7)$$

where

$$\begin{aligned} F_{y_{n_1}, y_{n_2}}(f) &= \frac{-(y_{n_1} + y_{n_2} - y_{n_1} y_{n_2})}{y_{n_1} y_{n_2}} \log(y_{n_1} + y_{n_2} - y_{n_1} y_{n_2}) \\ &+ \frac{(y_{n_1} + y_{n_2} - y_{n_1} y_{n_2})}{y_{n_1} y_{n_2}} \log(y_{n_1} + y_{n_2}) + \frac{y_{n_1}(1 - y_{n_2})}{y_{n_2}(y_{n_1} + y_{n_2})} \log(1 - y_{n_2}) \\ &+ \frac{y_{n_2}(1 - y_{n_1})}{y_{n_1}(y_{n_1} + y_{n_2})} \log(1 - y_{n_1}), \end{aligned}$$

is derived by use of the density of $F_{y_{n_1}, y_{n_2}}$ in §5. Because \widetilde{L}_1 and L_1 have the same asymptotic distribution and by (4.7), we get by letting $n_1 \wedge n_2 \rightarrow \infty$,

$$\widetilde{T}_N = v(f)^{-\frac{1}{2}} \left[-\frac{2 \log L_1}{N} - p \cdot F_{y_{n_1}, y_{n_2}}(f) - m(f) \right] \Rightarrow N(0, 1).$$

□

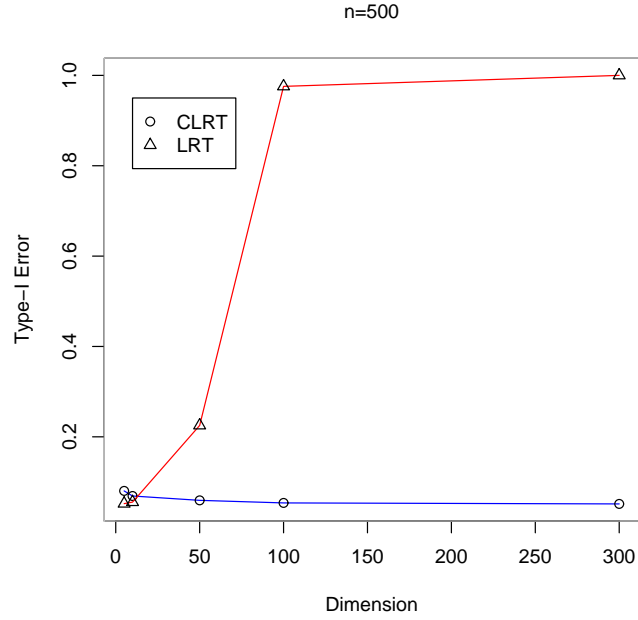


FIGURE 1. Realized sizes of the traditional LRT and the corrected LRT for different dimensions p with real Gaussian variables. 10 000 independent runs with 5% nominal level and sample size $n = 500$.

(y1, y2)=(0.05, 0.05)					
(p, n ₁ , n ₂)	CLRT			LRT	
	Size	Difference with 5%	Power	Size	Power
(5, 100, 100)	0.0770	0.0270	1	0.0582	1
(10, 200, 200)	0.0680	0.0180	1	0.0684	1
(20, 400, 400)	0.0593	0.0093	1	0.0872	1
(40, 800, 800)	0.0526	0.0026	1	0.1339	1
(80, 1600, 1600)	0.0501	0.0001	1	0.2687	1
(160, 3200, 3200)	0.0491	-0.0009	1	0.6488	1
(320, 6400, 6400)	0.0447	-0.0053	0.9671	1	1

(y1, y2)=(0.05, 0.1)					
(p, n ₁ , n ₂)	CLRT			LRT	
	Size	Difference with 5%	Power	Size	Power
(5, 100, 50)	0.0781	0.0281	0.9925	0.0640	0.9849
(10, 200, 100)	0.0617	0.0117	0.9847	0.0752	0.9904
(20, 400, 200)	0.0573	0.0073	0.9775	0.1104	0.9938
(40, 800, 400)	0.0561	0.0061	0.9765	0.2115	0.9975
(80, 1600, 800)	0.0521	0.0021	0.9702	0.4954	0.9998
(160, 3200, 1600)	0.0520	0.0020	0.9702	0.9433	1
(320, 6400, 3200)	0.0510	0.0010	1	0.9939	1

TABLE 2

Sizes and powers of the traditional LRT and the corrected LRT based on 10,000 independent replications using real Gaussian variables. Powers are estimated under the alternative $\Sigma_1 \Sigma_2^{-1} = \text{diag}(3, 1, 1, 1, \dots)$. Upper: $y_1 = y_2 = 0.05$. Bottom: $y_1 = 0.05, y_2 = 0.1$.

4.1. Simulation study II

For different values of (p, n_1, n_2) , we compute the realized sizes of the traditional LRT and the corrected LRT with 10,000 independent replications. The nominal test level is $\alpha = 0.05$ and we use real Gaussian variables. Results are summarized in Table 2 and Figure 2.

As we can see, when the dimension p increases, the traditional LRT leads to a dramatically high test size while the corrected LRT remains accurate. Furthermore, for moderate dimensions like $p = 20$ or 40 , the sizes of the traditional LRT are much higher than 5%, whereas the ones of corrected LRT are very close. By a closer look at the column showing the difference with 5%, we note that this difference rapidly decrease as p increases for the corrected test. Figure 2 gives a vivid sight of these comparisons between the traditional LRT and the corrected LRT in term of test sizes.

4.2. A pseudo-likelihood test for high-dimensional non-Gaussian data

As said in Introduction, previous related works as [Ledoit and Wolf \(2002\)](#), [Srivastava \(2005\)](#) or [Schott \(2007\)](#) all assume Gaussian variables. In contrast, Theorem 4.1 applies for general distributions having a fourth moment. For these non Gaussian data, we consider the corrected LRT as generalized pseudo-likelihood ratio test (or Gaussian LRT).

Moreover, the methods proposed by these authors all rely on an appropriate normalization of the trace of squared difference between two sample covariances following the idea of [Bai and Saranadasa \(1996\)](#). We believe that their method would strongly depend on the normality assumption (which was supported by simulation results below). On the other hand, based on general understanding, the LRT contains much higher information from data and its poor performance observed up to now is just caused by its large bias when dimension is large. Thus, from the intuitive understanding, we are confined ourselves to modify the LRT.

Let us develop in more details an example. Assume that \mathbf{x} follows a normalized t -distribution with 5 degree of freedom, that is $\mathbf{x} = \sqrt{\frac{3}{5}}t(5)$, \mathbf{x} and \mathbf{y} are i.i.d., hence $E\mathbf{x} = E\mathbf{y} = 0$, $E|\mathbf{x}|^2 = E|\mathbf{y}|^2 = 1$ and $E|\mathbf{x}|^4 = E|\mathbf{y}|^4 = 9$. We still employ the result in Theorem 4.1 for the test of equality between two covariance matrices, where

$$m_1(f) = \frac{1}{2} \left[\log \left(\frac{y_1 + y_2 - y_1 y_2}{y_1 + y_2} \right) - \frac{y_1}{y_1 + y_2} \log(1 - y_2) - \frac{y_2}{y_1 + y_2} \log(1 - y_1) \right. \\ \left. + \frac{6y_1^2 y_2}{(y_1 + y_2)^2} + \frac{6y_1 y_2^2}{(y_1 + y_2)^2} \right] \quad (4.8)$$

and

$$v_1(f) = -\frac{2y_2^2}{(y_1 + y_2)^2} \log(1 - y_1) - \frac{2y_1^2}{(y_1 + y_2)^2} \log(1 - y_2) - 2 \log \frac{y_1 + y_2}{y_1 + y_2 - y_1 y_2} \quad (4.9)$$

instead of $m(f)$ and $v(f)$ for real case, respectively. (4.8) and (4.9) are calculated in §5.

The following Table 3 summarizes a simulation study where we compare this corrected pseudo-LRT with the test proposed in Schott (2007). We use 1,000 independent replications with the above t -distributed variables. Again, the nominal test level is $\alpha = 0.05$. As we can see, the corrected pseudo-LRT performs correctly while Schott's test is no more valid here since the variables are not Gaussian.

(y1, y2)=(0.05, 0.1)		
(p, n1, n2)	CLRT Size	Schott's Size
(10,100, 200)	0.067	0.517
(20, 200, 400)	0.065	0.603
(40, 400, 800)	0.054	0.703
(80, 800, 1600)	0.048	0.764
(160, 1600, 3200)	0.045	0.826
(320, 3200, 6400)	0.051	0.854

TABLE 3

Sizes of the corrected pseudo-likelihood ration test and Schott's test for the case of $y_1 = 0.1$, $y_2 = 0.05$, based on 1,000 independent replications with normalized t -distributed variables with 5 degrees of freedom.

5. Proofs

Proof of (3.3)

By Theorem 2.1, for $g(x) = x - \log x - 1$, by using the variable change $x = 1 + y - 2\sqrt{y} \cos \theta$, $0 \leq \theta \leq \pi$, we have

$$\begin{aligned}
m(g) &= \frac{g(a(y)) + g(b(y))}{4} - \frac{1}{2\pi} \int_{a(y)}^{b(y)} \frac{g(x)}{\sqrt{4y - (x-1-y)^2}} dx \\
&= \frac{y - \log(1-y)}{2} - \frac{1}{2\pi} \int_0^\pi [1 + y - 2\sqrt{y} \cos \theta - \log(1 + y - 2\sqrt{y} \cos \theta) - 1] d\theta \\
&= \frac{y - \log(1-y)}{2} - \frac{1}{4\pi} \int_0^{2\pi} [y - 2\sqrt{y} \cos \theta - \log |1 - \sqrt{y} e^{i\theta}|^2] d\theta \\
&= -\frac{\log(1-y)}{2},
\end{aligned}$$

where $\int_0^{2\pi} \log |1 - \sqrt{y} e^{i\theta}|^2 d\theta = 0$ is calculated in Bai and Silverstein (2004).

Proof of (3.4)

For $g(x) = x - \log x - 1$, by Theorem 2.1, we have

$$v(g) = -\frac{1}{2\pi^2} \oint \oint \frac{g(z_1)g(z_2)}{(\underline{m}(z_1) - \underline{m}(z_2))^2} d\underline{m}(z_1) d\underline{m}(z_2)$$

and

$$\begin{aligned} g(z_1)g(z_2) &= z_1 z_2 - z_1 \log z_2 - z_2 \log z_1 + \log z_1 \log z_2 \\ &\quad - z_1 + \log z_1 - z_2 + \log z_2 + 1. \end{aligned}$$

It is easy to see that $v(\mathbf{1}, \mathbf{1}) = 0$, where $\mathbf{1}$ means constant function equals to 1. For Stieltjes transform of F^y , the following equation is given in [Bai and Silverstein \(2004\)](#), for $z \in \mathbb{C}^+$,

$$z = -\frac{1}{\underline{m}(z)} + \frac{y}{1 + \underline{m}(z)}. \quad (5.1)$$

Let $m_i = \underline{m}(z_i)$, $i = 1, 2$. For fixed m_2 , we have on a contour enclosed 1, $(y-1)^{-1}$ and -1, but not 0,

$$\begin{aligned} \oint \frac{\log(z(m_1))}{(m_1 - m_2)^2} dm_1 &= \oint \frac{\frac{1}{m_1^2} - \frac{y}{(1+m_1)^2}}{-\frac{1}{m_1} + \frac{y}{1+m_1}} \frac{1}{(m_1 - m_2)} dm_1 \\ &= \oint \frac{(1+m_1)^2 - ym_1^2}{ym_1(m_1 - m_2)} \left(\frac{-1}{m_1 + 1} + \frac{1}{m_1 - \frac{1}{y-1}} \right) dm_1 \\ &= 2\pi i \cdot \left(\frac{1}{m_2 + 1} - \frac{1}{m_2 - \frac{1}{y-1}} \right). \end{aligned}$$

and

$$\begin{aligned} &\oint \frac{-\frac{1}{m_1} + \frac{y}{1+m_1}}{(m_1 - m_2)^2} dm_1 \\ &= y \oint \left(\frac{1}{1+m_1} + \frac{1-y}{y} \right) \cdot [1 - (1+m_1)]^{-1} \cdot (m_2 + 1)^{-2} \cdot \left(1 - \frac{m_1 + 1}{m_2 + 1}\right)^{-2} dm_1 \\ &= y \oint \left(\frac{1}{1+m_1} + \frac{1-y}{y} \right) \cdot \sum_{j=0}^{\infty} (1+m_1)^j (m_2 + 1)^{-2} \sum_{\ell=1}^{\infty} \ell \left(\frac{m_1 + 1}{m_2 + 1} \right)^{\ell-1} dm_1 \\ &= 2\pi i \cdot \frac{y}{(m_2 + 1)^2}. \end{aligned}$$

Then we also get $v(-z_1 + \log z_1, \mathbf{1}) = 0$. Similarly, $v(\mathbf{1}, -z_2 + \log z_2) = 0$. Furthermore,

$$v(z_1, z_2) = \frac{y^2}{\pi i} \oint \frac{1}{(m_2 + 1)^2} \left(\frac{1}{1+m_2} + \frac{1-y}{y} \right) \sum_{j=0}^{\infty} (1+m_2)^j dm_2 = 2y,$$

and

$$\begin{aligned} v(z_1, \log z_2) &= \frac{y}{\pi i} \oint \left(\frac{1}{m_2 + 1} - \frac{1}{m_2 - 1/(y-1)} \right) \left(\frac{1}{1+m_2} + \frac{1-y}{y} \right) \cdot [1 - (1+m_2)]^{-1} dm_2 \\ &= \frac{y}{\pi i} \oint \left(\frac{1}{m_2 + 1} - \frac{1}{m_2 - 1/(y-1)} \right) \left(\frac{1}{1+m_2} + \frac{1-y}{y} \right) \sum_{j=0}^{\infty} (1+m_2)^j dm_2 \\ &= 2y. \end{aligned}$$

By a computation in [Bai and Silverstein \(2004\)](#), we know that $v(\log z_1, \log z_2) = -2 \log(1-y)$.

Finally, we obtain

$$v(g) = v(z_1, z_2) + v(\log z_1, \log z_2) - 2v(z_1, \log z_2)$$

$$\begin{aligned}
& +v(-z_1 + \log z_1, \mathbf{1}) + v(\mathbf{1}, -z_2 + \log z_2) + v(\mathbf{1}, \mathbf{1}) \\
& = -2\log(1-y) - 2y.
\end{aligned}$$

Proof of (3.6)

Since F^{y_n} is the Marčenko-Pastur law of index y_n , by using the variable change $x = 1 + y_n - 2\sqrt{y_n}\cos\theta$, $0 \leq \theta \leq \pi$ we have

$$\begin{aligned}
F^{y_n}(g) &= \int_{a(y_n)}^{b(y_n)} \frac{x - \log x - 1}{2\pi xy_n} \sqrt{(b(y_n) - x)(x - a(y_n))} dx \\
&= \frac{1}{2\pi y_n} \int_0^\pi \left[1 - \frac{\log(1 + y_n - 2\sqrt{y_n}\cos\theta) + 1}{1 + y_n - 2\sqrt{y_n}\cos\theta} \right] 4y_n \sin^2\theta d\theta \\
&= \frac{1}{2\pi} \int_0^{2\pi} \left[2\sin^2\theta - \frac{2\sin^2\theta}{1 + y_n - 2\sqrt{y_n}\cos\theta} (\log|1 - \sqrt{y_n}e^{i\theta}|^2 - 1) \right] d\theta \\
&= 1 - \frac{y_n - 1}{y_n} \log(1 - y_n),
\end{aligned}$$

where

$$\frac{1}{2\pi} \int_0^{2\pi} \frac{2\sin^2\theta}{1 + y_n - 2\sqrt{y_n}\cos\theta} \log|1 - \sqrt{y_n}e^{i\theta}|^2 d\theta = \frac{y_n - 1}{y_n} \log(1 - y_n) - 1$$

is calculated in [Bai and Silverstein \(2004\)](#).

Proof of Lemma 2.1

We use the variable change $x = (1 - y_2)^{-2}(1 + h^2 - 2h\cos\theta)$, where $h = \sqrt{y_1 + y_2 - y_1 y_2}$. When c, d satisfy $c^2 + d^2 = a(1 - y_2)^2 + b(1 + h^2)$, $cd = bh$, $0 < d < c$, we have

$$f(z(\xi)) = \log(a + bz(\xi)) = \log\left(\frac{|c + d\xi|^2}{(1 - y_2)^2}\right).$$

Similarly,

$$g(z(\xi)) = \log(\alpha + \beta z(\xi)) = \log\left(\frac{|\gamma + \eta\xi|^2}{(1 - y_2)^2}\right).$$

Let

$$\tilde{f}(z(\xi)) = \log\left(\frac{(c + d\xi)^2}{(1 - y_2)^2}\right) \quad \text{and} \quad \tilde{g}(z(\xi)) = \log\left(\frac{(\gamma + \eta\xi)^2}{(1 - y_2)^2}\right).$$

Note that $f(z(\xi)) = \Re(\tilde{f}(z(\xi)))$ and $g(z(\xi)) = \Re(\tilde{g}(z(\xi)))$. By Theorem 2.2, we have

$$\begin{aligned}
m(f) &= \frac{1}{4\pi i} \oint_{|\xi|=1} f(z(\xi)) \left[\frac{1}{\xi - \frac{1}{r}} + \frac{1}{\xi + \frac{1}{r}} - \frac{2}{\xi + \frac{y_2}{hr}} \right] d\xi \\
&= \frac{1}{4\pi} \int_0^{2\pi} f(z(e^{i\theta})) \left[\frac{1}{e^{i\theta} - \frac{1}{r}} + \frac{1}{e^{i\theta} + \frac{1}{r}} - \frac{2}{e^{i\theta} + \frac{y_2}{hr}} \right] e^{i\theta} d\theta \\
&= \frac{1}{4\pi} \int_0^{2\pi} f(z(e^{i\theta})) \left[\frac{1}{e^{-i\theta} - \frac{1}{r}} + \frac{1}{e^{-i\theta} + \frac{1}{r}} - \frac{2}{e^{-i\theta} + \frac{y_2}{hr}} \right] e^{-i\theta} d\theta
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{8\pi} \int_0^{2\pi} f(z(e^{i\theta})) \left\{ \left[\frac{1}{e^{i\theta} - \frac{1}{r}} + \frac{1}{e^{i\theta} + \frac{1}{r}} - \frac{2}{e^{i\theta} + \frac{y_2}{hr}} \right] e^{i\theta} + \right. \\
&\quad \left. \left[\frac{1}{e^{-i\theta} - \frac{1}{r}} + \frac{1}{e^{-i\theta} + \frac{1}{r}} - \frac{2}{e^{-i\theta} + \frac{y_2}{hr}} \right] e^{-i\theta} \right\} d\theta \\
&= \frac{1}{8\pi} \Re \left\{ \int_0^{2\pi} \tilde{f}(z(e^{i\theta})) \left[\left(\frac{1}{e^{i\theta} - \frac{1}{r}} + \frac{1}{e^{i\theta} + \frac{1}{r}} - \frac{2}{e^{i\theta} + \frac{y_2}{hr}} \right) e^{i\theta} + \right. \right. \\
&\quad \left. \left. \left(\frac{r}{r - e^{i\theta}} + \frac{r}{r + e^{i\theta}} - \frac{2hr}{y_2 e^{i\theta} + hr} \right) \right] d\theta \right\} \\
&= \Re \left\{ \frac{1}{8\pi i} \oint_{|\xi|=1} \tilde{f}(z(\xi)) \left[\left(\frac{1}{\xi - \frac{1}{r}} + \frac{1}{\xi + \frac{1}{r}} - \frac{2}{\xi + \frac{y_2}{hr}} \right) \right. \right. \\
&\quad \left. \left. + \left(\frac{r}{r - \xi} + \frac{r}{r + \xi} - \frac{2hr}{y_2 \xi + hr} \right) \xi^{-1} \right] d\xi \right\} \\
&= \frac{1}{4} \left(\tilde{f}\left(z\left(\frac{1}{r}\right)\right) + \tilde{f}\left(z\left(-\frac{1}{r}\right)\right) - 2\tilde{f}\left(z\left(-\frac{y_2}{hr}\right)\right) \right) \\
&\rightarrow r \downarrow 1 \frac{1}{4} \left[\tilde{f}(z(1)) + \tilde{f}(z(-1)) - 2\tilde{f}\left(z\left(-\frac{y_2}{h}\right)\right) \right] \\
&= \frac{1}{2} \log \frac{(c^2 - d^2)h^2}{(ch - y_2 d)^2}.
\end{aligned}$$

Let $m_j = -\frac{1+hr_j\xi_j}{1-y_2}$, where $|\xi_j| = 1, j = 1, 2, r_2 \downarrow r_1$, and $r_1 \downarrow 1$. By Theorem 2.2, we have

$$v(f, g) = -\frac{1}{2\pi^2} \oint_{|\xi_2|=1} \left\{ \oint_{|\xi_1|=1} \frac{f(z(r_1\xi_1))}{(r_2\xi_2 - r_1\xi_1)^2} \cdot r_1 r_2 d\xi_1 \right\} g(z(r_2\xi_2)) d\xi_2.$$

When $r_1 \downarrow 1$, $-\frac{d}{cr_1}$ and 0 are poles. We can then choose r_1 so that $-\frac{c}{dr_1}$ is not a pole. Then we get

$$\begin{aligned}
&\oint_{|\xi_1|=1} \frac{\log(a + bz(r_1\xi_1))}{(r_2\xi_2 - r_1\xi_1)^2} \cdot r_1 r_2 d\xi_1 \\
&= \oint_{|\xi_1|=1} \frac{(\log(a + bz(r_1\xi_1)))'}{r_1\xi_1 - r_2\xi_2} \cdot r_2 d\xi_1 \\
&= \oint_{|\xi_1|=1} \left[\frac{bhr_1\xi_1}{(r_1\xi_1 - r_2\xi_2)(c + dr_1\xi_1)c} \cdot \frac{1}{\xi_1 + \frac{d}{cr_1}} \right. \\
&\quad \left. - \frac{bhr_1^{-1}}{(r_1\xi_1 - r_2\xi_2)(c + dr_1\xi_1)c} \cdot \frac{1}{(\xi_1 + \frac{d}{cr_1})\xi_1} \cdot r_2 \right] d\xi_1 \\
&= 2\pi i \left(\frac{bhd^{-1}c^{-1}}{\xi_2} - \frac{bhd^{-1}r_2}{d + cr_2\xi_2} \right).
\end{aligned}$$

So,

$$v(f, g) = -\frac{i}{\pi} \oint_{|\xi_2|=1} \left(\frac{bhd^{-1}c^{-1}}{\xi_2} - \frac{bhd^{-1}r_2}{d + cr_2\xi_2} \right) \log(\alpha + \beta z(r_2\xi_2)) d\xi_2.$$

Since the function $g(x) = \log(\alpha + \beta x)$ is analytic, when $r_2 > 1$ but sufficiently close to 1, we have

$$|g(z(r_2\xi_2)) - g(z(\xi_2))| \leq K(r - 1),$$

for some constant K . Thus we have

$$\left| \oint_{|\xi|_2=1} [g(z(r_2\xi_2)) - g(z(\xi_2))] \left(\frac{bhd^{-1}c^{-1}}{\xi_2} - \frac{bhd^{-1}r_2}{d + cr_2\xi_2} \right) d\xi_2 \right| \rightarrow 0 \quad \text{as } r_2 \downarrow 1,$$

where the estimations are done according to $|\arg(\xi_2)|$ or $|\arg(\xi_2) - \pi| \leq \sqrt{r_2 - 1}$ or not. Thus,

$$v(f, g) = -\frac{i}{\pi} \oint_{|\xi_2|=1} g(z(\xi_2)) \left(\frac{bhd^{-1}c^{-1}}{\xi_2} - \frac{bhd^{-1}r_2}{d + cr_2\xi_2} \right) d\xi_2 + R(r_2)$$

where $R(r_2) \rightarrow 0$, as $r_2 \downarrow 1$. Because $g(z(\xi_2)) = \log \left(\frac{|\gamma + \eta\xi_2|^2}{(1 - y_2)^2} \right)$, for γ, η satisfying $\gamma^2 + \eta^2 =$

$\alpha(1 - y_2)^2 + \beta(1 + h^2)$, $\gamma\eta = \beta h$, $0 < \eta < \gamma$, and if $\tilde{g}(z(\xi_2)) = \log \left(\frac{(\gamma + \eta\xi_2)^2}{(1 - y_2)^2} \right)$, we have $g(z(\xi_2)) = \Re(\tilde{g}(z(\xi_2)))$. Therefore,

$$\begin{aligned} v(f, g) &= -\frac{i}{\pi} \oint_{|\xi|_2=1} g(z(\xi_2)) \left(\frac{bhd^{-1}c^{-1}}{\xi_2} - \frac{bhd^{-1}r_2}{d + cr_2\xi_2} \right) d\xi_2 \\ &= \frac{1}{\pi} \int_0^{2\pi} g(z(e^{i\theta})) \left(\frac{bhd^{-1}c^{-1}}{e^{i\theta}} - \frac{bhd^{-1}r_2}{d + cr_2e^{i\theta}} \right) e^{i\theta} d\theta \\ &= \frac{\theta \rightarrow 2\pi - \theta}{\pi} \int_0^{2\pi} g(z(e^{i\theta})) \left(\frac{bhd^{-1}c^{-1}}{e^{-i\theta}} - \frac{bhd^{-1}r_2}{d + cr_2e^{-i\theta}} \right) e^{-i\theta} d\theta \\ &= \frac{1}{2\pi} \int_0^{2\pi} g(z(e^{i\theta})) \left[\left(\frac{bhd^{-1}c^{-1}}{e^{i\theta}} - \frac{bhd^{-1}r_2}{d + cr_2e^{i\theta}} \right) e^{i\theta} + bhd^{-1}c^{-1} - \frac{bhd^{-1}r_2}{de^{i\theta} + cr_2} \right] d\theta \\ &= \frac{1}{2\pi} \Re \left\{ \int_0^{2\pi} \tilde{g}(z(e^{i\theta})) \left[\left(\frac{bhd^{-1}c^{-1}}{e^{i\theta}} - \frac{bhd^{-1}r_2}{d + cr_2e^{i\theta}} \right) e^{i\theta} + bhd^{-1}c^{-1} - \frac{bhd^{-1}r_2}{de^{i\theta} + cr_2} \right] d\theta \right\} \\ &= \Re \left\{ \frac{1}{2\pi i} \oint_{|\xi|_2=1} \tilde{g}(z(\xi_2)) \left[\left(\frac{bhd^{-1}c^{-1}}{\xi_2} - \frac{bhd^{-1}r_2}{d + cr_2\xi_2} \right) + \left(bhd^{-1}c^{-1} - \frac{bhd^{-1}r_2}{d\xi_2 + cr_2} \right) \xi_2^{-1} \right] d\xi_2 \right\} \\ &= bhd^{-1}c^{-1} \left[\tilde{g}(z(0)) - \tilde{g}(z(-\frac{d}{cr_2})) \right] \\ &\rightarrow bhd^{-1}c^{-1} \left[\tilde{g}(z(0)) - \tilde{g}(z(-\frac{d}{c})) \right] \\ &= 2bhd^{-1}c^{-1} \log \frac{c\gamma}{c\gamma - d\eta}. \end{aligned}$$

Proof of (4.5) and (4.6)

Because ξ and η are Gaussian variables, for real case, $\beta = E|\xi|^4 - 3 = 0$, then (2.10), (2.11) and (2.13) are all 0. Consider (2.9) and (2.12), as $y_{n_k} \rightarrow y_k$, $k = 1, 2$, by the computations done in the proof of Lemma 2.1, we see that termes tending to zero could be neglected in the considered contour integrals. Hence we can put $y_{n_k} = y_k$, $k = 1, 2$ and use

$$f(x) = \log(y_1 + y_2x) - \frac{y_2}{y_1 + y_2} \log x - \log(y_1 + y_2)$$

instead of $f(x) = \log(y_{n_1} + y_{n_2}x) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log x - \log(y_{n_1} + y_{n_2})$. Consider the variable change $x = (1-y_2)^{-2}(1+h^2-2h \cos \theta)$, where $z(\xi) = (1-y_2)^{-2} [1 + h^2 + 2h\mathcal{R}(\xi)]$, $h = \sqrt{y_1 + y_2 - y_1 y_2}$.

As

$$\begin{aligned} \log(y_{n_1} + y_{n_2}z(\xi)) &= \log \left(\frac{|h + y_2\xi|^2}{(1-y_2)^2} \right), \\ \log(z(\xi)) &= \log \left(\frac{|1 + h\xi|^2}{(1-y_2)^2} \right), \end{aligned}$$

we have by Lemma 2.1,

$$\begin{aligned} m(f) &= \frac{1}{2} \left[\log \frac{(h^2 - y_2^2)h^2}{(h^2 - y_2^2)^2} - \frac{y_2}{y_1 + y_2} \log \frac{(1-h^2)h^2}{(h - y_2h)^2} \right] \\ &= \frac{1}{2} \left[\log \left(\frac{y_1 + y_2 - y_1 y_2}{y_1 + y_2} \right) - \frac{y_1}{y_1 + y_2} \log(1 - y_2) - \frac{y_2}{y_1 + y_2} \log(1 - y_1) \right], \end{aligned}$$

and

$$\begin{aligned} v(f) &= v(\log(y_{n_1} + y_{n_2}x)) + \frac{y_2^2}{(y_1 + y_2)^2} v(\log x) - \frac{2y_2}{y_1 + y_2} v(\log x, \log(y_{n_1} + y_{n_2}x)) \\ &= 2 \log \frac{h^2}{h^2 - y_2^2} + 2 \frac{y_2^2}{(y_1 + y_2)^2} \log \frac{1}{1 - h^2} - \frac{4y_2}{y_1 + y_2} \log \frac{1}{1 - y_2} \\ &= -\frac{2y_2^2}{(y_1 + y_2)^2} \log(1 - y_1) - \frac{2y_1^2}{(y_1 + y_2)^2} \log(1 - y_2) - 2 \log \frac{y_1 + y_2}{y_1 + y_2 - y_1 y_2}. \end{aligned}$$

Proof of $F_{y_{n_1}, y_{n_2}}(f)$

By (4.3) and the density of $F_{y_{n_1}, y_{n_2}}(f)$ (the limiting distribution in (2.7) but with y_{n_k} in place of $y_k, k = 1, 2$), where $h_n = \sqrt{y_{n_1} + y_{n_2} - y_{n_1} y_{n_2}}$, $a_n = (1 - y_{n_2})^{-2} (1 - \sqrt{y_{n_1} + y_{n_2} - y_{n_1} y_{n_2}})^2$ and $b_n = (1 - y_{n_2})^{-2} (1 + \sqrt{y_{n_1} + y_{n_2} - y_{n_1} y_{n_2}})^2$. Using the substitution $x = (1 - y_{n_2})^{-2} (1 + h_n^2 - 2h_n \cos \theta)$, $0 < \theta < \pi$, we have

$$\begin{aligned} \sqrt{(b_n - x)(x - a_n)} &= \frac{2h_n \sin \theta}{(1 - y_{n_2})^2}, & dx &= \frac{2h_n \sin \theta d\theta}{(1 - y_{n_2})^2}; \\ x &= \frac{|1 - h_n e^{i\theta}|^2}{(1 - y_{n_2})^2}, & y_{n_1} + y_{n_2}x &= \frac{|h_n - y_{n_2} e^{i\theta}|^2}{(1 - y_{n_2})^2}. \end{aligned}$$

Therefore,

$$\begin{aligned} &F_{y_{n_1}, y_{n_2}}(f) \\ &= \int_{a_n}^{b_n} f(x) \frac{(1 - y_{n_2}) \sqrt{(b_n - x)(x - a_n)}}{2\pi x (y_{n_1} + y_{n_2}x)} dx \\ &= (1 - y_{n_2}) \int_{a_n}^{b_n} \left[\log(y_{n_1} + y_{n_2}x) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log x \right] \frac{\sqrt{(b_n - x)(x - a_n)}}{2\pi x (y_{n_1} + y_{n_2}x)} dx \\ &\quad - \log(y_{n_1} + y_{n_2}) \end{aligned}$$

$$\begin{aligned}
&= \frac{2(1-y_{n_2})}{\pi} \int_0^\pi \left[\log \frac{|h_n - y_{n_2} e^{i\theta}|^2}{(1-y_{n_2})^2} - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log \frac{|1 - h_n e^{i\theta}|^2}{(1-y_{n_2})^2} \right] \\
&\quad \cdot \frac{h_n^2 \sin^2 \theta}{|1 - h_n e^{i\theta}|^2 |h_n - y_{n_2} e^{i\theta}|^2} d\theta - \log(y_{n_1} + y_{n_2}) \\
&= \frac{2(1-y_{n_2})}{\pi} \int_0^\pi \left[\log |h_n - y_{n_2} e^{i\theta}|^2 - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log |1 - h_n e^{i\theta}|^2 \right] \\
&\quad \cdot \frac{h_n^2 \sin^2 \theta}{|1 - h_n e^{i\theta}|^2 |h_n - y_{n_2} e^{i\theta}|^2} d\theta - 2 \left(1 - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \right) \log(1 - y_{n_2}) - \log(y_{n_1} + y_{n_2}) \\
&= \Re \left\{ \frac{2(1-y_{n_2})}{\pi} \int_0^{2\pi} \left[\log(h_n - y_{n_2} e^{i\theta}) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log(1 - h_n e^{i\theta}) \right] \right. \\
&\quad \left. \cdot \frac{h_n^2 \sin^2 \theta}{|1 - h_n e^{i\theta}|^2 |h_n - y_{n_2} e^{i\theta}|^2} d\theta \right\} - \frac{2y_{n_1}}{y_{n_1} + y_{n_2}} \log(1 - y_{n_2}) - \log(y_{n_1} + y_{n_2}) \\
&= \Re \left\{ \frac{-(1-y_{n_2})}{2\pi i} \oint_{|z|=1} \left[\log(h_n - y_{n_2} z) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log(1 - h_n z) \right] \right. \\
&\quad \left. \cdot \frac{h_n^2 (z - z^{-1})^2}{z |1 - h_n z|^2 |h_n - y_{n_2} z|^2} dz \right\} - \frac{2y_{n_1}}{y_{n_1} + y_{n_2}} \log(1 - y_{n_2}) - \log(y_{n_1} + y_{n_2}) \\
&= \Re \left\{ \frac{y_{n_2} - 1}{y_{n_2}} \cdot \frac{1}{2\pi i} \oint_{|z|=1} \left[\log(h_n - y_{n_2} z) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log(1 - h_n z) \right] \right. \\
&\quad \left. \cdot \frac{(z^2 - 1)^2}{z(z - h_n)(z - \frac{1}{h_n})(z - \frac{y_{n_2}}{h_n})(z - \frac{h_n}{y_{n_2}})} dz \right\} - \frac{2y_{n_1}}{y_{n_1} + y_{n_2}} \log(1 - y_{n_2}) - \log(y_{n_1} + y_{n_2}).
\end{aligned}$$

There are three poles inside the unit circle: 0, $h_n, y_{n_2}/h_n$. Their corresponding residues are

$$\begin{aligned}
R(0) &= \frac{y_{n_2} - 1}{y_{n_2}} \log(h_n), \\
R(h_n) &= \frac{(h_n^2 - 1)}{(h_n^2 - y_{n_2})} \left[\log(h_n) + \log(1 - y_{n_2}) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log(1 - h_n^2) \right], \\
R\left(\frac{y_{n_2}}{h_n}\right) &= \frac{(y_{n_2}^2 - h_n^2)}{y_{n_2}(y_{n_2} - h_n^2)} \left[\log(h_n^2 - y_{n_2}^2) - \log(h_n) - \frac{y_{n_2}}{y_{n_1} + y_{n_2}} \log(1 - y_{n_2}) \right].
\end{aligned}$$

Therefore,

$$\begin{aligned}
F^{y_{n_1}, y_{n_2}}(f) &= R(0) + R(h_n) + R\left(\frac{y_{n_2}}{h_n}\right) - \frac{2y_{n_1}}{y_{n_1} + y_{n_2}} \log(1 - y_{n_2}) - \log(y_{n_1} + y_{n_2}) \\
&= \frac{-(y_{n_1} + y_{n_2} - y_{n_1} y_{n_2})}{y_{n_1} y_{n_2}} \log(y_{n_1} + y_{n_2} - y_{n_1} y_{n_2}) \\
&\quad + \frac{(y_{n_1} + y_{n_2} - y_{n_1} y_{n_2})}{y_{n_1} y_{n_2}} \log(y_{n_1} + y_{n_2}) + \frac{y_{n_1}(1 - y_{n_2})}{y_{n_2}(y_{n_1} + y_{n_2})} \log(1 - y_{n_2})
\end{aligned}$$

$$+ \frac{y_{n_2}(1 - y_{n_1})}{y_{n_1}(y_{n_1} + y_{n_2})} \log(1 - y_{n_1}).$$

Proof of (4.8) and (4.9)

Because \mathbf{x} and \mathbf{y} are random variables from normalized t -distribution with 5 degree of freedom, \mathbf{x} and \mathbf{y} are *i.i.d.*, $E\mathbf{x} = E\mathbf{y} = 0$, $E|\mathbf{x}|^2 = E|\mathbf{y}|^2 = 1$ and $E|\mathbf{x}|^4 = E|\mathbf{y}|^4 = 9$. For real case, $\beta = E|\xi|^4 - 3 = 6$, (2.9) and (2.12) items are the same to the Gaussian variables. Consider the items (2.10), (2.11) and (2.13). As the same explanation in Proof of (4.5) and (4.6), we use $f(x) = \log(y_1 + y_2x) - \frac{y_2}{y_1 + y_2} \log x - \log(y_1 + y_2)$ instead.

For (2.10), we have

$$\begin{aligned} & \frac{\beta \cdot y_1(1 - y_2)^2}{2\pi i \cdot h^2} \oint_{|\xi|=1} \left[\log \frac{|h + y_2\xi|^2}{(1 - y_2)^2} - \frac{y_2}{y_1 + y_2} \log \frac{|1 + h\xi|^2}{(1 - y_2)^2} - \log(y_1 + y_2) \right] \\ & \quad \cdot \frac{1}{(\xi + \frac{y_2}{hr})^3} d\xi \\ &= \frac{\beta \cdot y_1(1 - y_2)^2}{2\pi i \cdot h^2} \oint_{|\xi|=1} 2\mathcal{R} \left\{ \log(h + y_2\xi) - \frac{y_2}{y_1 + y_2} \log(1 + h\xi) \right\} \cdot \frac{1}{(\xi + \frac{y_2}{hr})^3} d\xi \\ &= \frac{\beta \cdot y_1(1 - y_2)^2}{2\pi i \cdot h^2} \oint_{|\xi|=1} \left\{ \log(h + y_2\xi) + \log(h + y_2\bar{\xi}) \right. \\ & \quad \left. - \frac{y_2}{y_1 + y_2} [\log(1 + h\xi) + \log(1 + h\bar{\xi})] \right\} \cdot \frac{1}{(\xi + \frac{y_2}{hr})^3} d\xi \\ &= \frac{\beta \cdot y_1(1 - y_2)^2}{2\pi i \cdot h^2} \oint_{|\xi|=1} \left[\log(h + y_2\xi) - \frac{y_2}{y_1 + y_2} \log(1 + h\xi) \right] \\ & \quad \left[\frac{1}{(\xi + \frac{y_2}{hr})^3} + \frac{\left(\frac{hr}{y_2}\right)^3 \xi}{(\frac{hr}{y_2} + \xi)^3} \right] d\xi \\ &= \frac{\beta \cdot y_1(1 - y_2)^2}{2\pi i \cdot h^2} \cdot 2\pi \cdot \frac{1}{2} \left[\log(h + y_2\xi) - \frac{y_2}{y_1 + y_2} \log(1 + h\xi) \right]'' \Big|_{\xi = -\frac{y_2}{hr}} \\ &= \frac{\beta \cdot y_1(1 - y_2)^2}{2h^2} \left[-\frac{y_2^2}{(h + y_2\xi)^2} + \frac{y_2}{y_1 + y_2} \frac{h^2}{(1 + h\xi)^2} \right] \Big|_{\xi = -\frac{y_2}{hr}} \\ &= \frac{\beta y_1^2 y_2}{2(y_1 + y_2)^2}. \end{aligned}$$

For (2.11), we have

$$\begin{aligned} & \frac{\beta \cdot (1 - y_2)}{4\pi i} \oint_{|\xi|=1} \left[\log \frac{|h + y_2\xi|^2}{(1 - y_2)^2} - \frac{y_2}{y_1 + y_2} \log \frac{|1 + h\xi|^2}{(1 - y_2)^2} - \log(y_1 + y_2) \right] \\ & \quad \cdot \frac{\xi^2 - \frac{y_2^2}{h^2 r^2}}{(\xi + \frac{y_2}{hr})^2} \left[\frac{1}{\xi - \frac{\sqrt{y_2}}{hr}} + \frac{1}{\xi + \frac{\sqrt{y_2}}{hr}} - \frac{2}{\xi + \frac{y_2}{hr}} \right] d\xi \\ &= \frac{\beta \cdot (1 - y_2) y_2}{2\pi i \cdot h} \oint_{|\xi|=1} \left[\log \frac{|h + y_2\xi|^2}{(1 - y_2)^2} - \frac{y_2}{y_1 + y_2} \log \frac{|1 + h\xi|^2}{(1 - y_2)^2} - \log(y_1 + y_2) \right] \\ & \quad \cdot \left[\frac{\xi + \frac{1}{hr}}{(\xi + \frac{y_2}{hr})^3} \right] d\xi \end{aligned}$$

$$\begin{aligned}
&= \frac{\beta \cdot (1-y_2)y_2}{2\pi i \cdot h} \oint_{|\xi|=1} 2\mathcal{R} \left\{ \log(h+y_2\xi) - \frac{y_2}{y_1+y_2} \log(1+h\xi) \right\} \cdot \left[\frac{\xi + \frac{1}{hr}}{(\xi + \frac{y_2}{hr})^3} \right] d\xi \\
&= \frac{\beta \cdot (1-y_2)y_2}{2\pi i \cdot h} \oint_{|\xi|=1} \left[\log(h+y_2\xi) + \log(h+y_2\bar{\xi}) \right. \\
&\quad \left. - \frac{y_2}{y_1+y_2} (\log(1+h\xi) + \log(1+h\bar{\xi})) \right] \cdot \left[\frac{\xi + \frac{1}{hr}}{(\xi + \frac{y_2}{hr})^3} \right] d\xi \\
&= \frac{\beta \cdot (1-y_2)y_2}{2\pi i \cdot h} \oint_{|\xi|=1} \left[\log(h+y_2\xi) - \frac{y_2}{y_1+y_2} \log(1+h\xi) \right] \\
&\quad \cdot \left[\frac{\xi + \frac{1}{hr}}{(\xi + \frac{y_2}{hr})^3} + \frac{\frac{h^2 r^2}{y_2^3} (\xi + hr)}{(\xi + \frac{hr}{y_2})^3} \right] d\xi \\
&= \frac{\beta \cdot (1-y_2)y_2}{2\pi i \cdot h} 2\pi i \cdot \frac{1}{2} \left(\left[\log(h+y_2\xi) - \frac{y_2}{y_1+y_2} \log(1+h\xi) \right] \cdot \left(\xi + \frac{1}{hr} \right) \right)' \bigg|_{\xi = -\frac{y_2}{hr}} \\
&= \frac{\beta y_1 y_2^2}{2(y_1+y_2)^2}.
\end{aligned}$$

Therefore,

$$\begin{aligned}
m_1(f) &= \frac{1}{2} \left[\log \left(\frac{y_1+y_2-y_1 y_2}{y_1+y_2} \right) - \frac{y_1}{y_1+y_2} \log(1-y_2) - \frac{y_2}{y_1+y_2} \log(1-y_1) \right. \\
&\quad \left. + \frac{6y_1^2 y_2}{(y_1+y_2)^2} + \frac{6y_1 y_2^2}{(y_1+y_2)^2} \right].
\end{aligned}$$

For covariance, we have

$$\begin{aligned}
&\oint_{|\xi|=1} \frac{f \left(\frac{1+h^2+2h\mathcal{R}(\xi)}{(1-y_2)^2} \right)}{(\xi + \frac{y_2}{hr})^2} d\xi \\
&= \oint_{|\xi|=1} \left[\log \frac{|h+y_2\xi|^2}{(1-y_2)^2} - \frac{y_2}{y_1+y_2} \log \frac{|1+h\xi|^2}{(1-y_2)^2} - \log(y_1+y_2) \right] \cdot \frac{1}{(\xi + \frac{y_2}{hr})^2} d\xi \\
&= \oint_{|\xi|=1} 2\mathcal{R} \left\{ \log(h+y_2\xi) - \frac{y_2}{y_1+y_2} \log(1+h\xi) \right\} \cdot \frac{1}{(\xi + \frac{y_2}{hr})^2} d\xi \\
&= \oint_{|\xi|=1} \left\{ \log(h+y_2\xi) + \log(h+y_2\bar{\xi}) \right. \\
&\quad \left. - \frac{y_2}{y_1+y_2} [\log(1+h\xi) + \log(1+h\bar{\xi})] \right\} \cdot \frac{1}{(\xi + \frac{y_2}{hr})^2} d\xi \\
&= \oint_{|\xi|=1} \left[\log(h+y_2\xi) - \frac{y_2}{y_1+y_2} \log(1+h\xi) \right] \left[\frac{1}{(\xi + \frac{y_2}{hr})^2} + \frac{\left(\frac{hr}{y_2} \right)^2}{(\xi + \frac{hr}{y_2})^2} \right] d\xi \\
&= 2\pi i \cdot \left[\log(h+y_2\xi) - \frac{y_2}{y_1+y_2} \log(1+h\xi) \right]' \bigg|_{\xi = -\frac{y_2}{hr}} \\
&= \pi i \cdot \left[\frac{y_2}{h+y_2\xi} + \frac{y_2}{y_1+y_2} \frac{h}{1+h\xi} \right] \bigg|_{\xi = -\frac{y_2}{hr}} \\
&= 0.
\end{aligned}$$

So, (2.13) becomes,

$$-\frac{\beta \cdot (y_1 + y_2)(1 - y_2)^2}{4\pi^2 h^2} \oint_{|\xi_1|=1} \frac{f\left(\frac{1+h^2+2h\mathcal{R}(\xi_1)}{(1-y_2)^2}\right)}{(\xi_1 + \frac{y_2}{hr_1})^2} d\xi_1 \oint_{|\xi_2|=1} \frac{f\left(\frac{1+h^2+2h\mathcal{R}(\xi_2)}{(1-y_2)^2}\right)}{(\xi_2 + \frac{y_2}{hr_2})^2} d\xi_2 = 0$$

Finally,

$$v_1(f) = -\frac{2y_2^2}{(y_1 + y_2)^2} \log(1 - y_1) - \frac{2y_1^2}{(y_1 + y_2)^2} \log(1 - y_2) - 2 \log \frac{y_1 + y_2}{y_1 + y_2 - y_1 y_2}.$$

References

- Anderson, T. W. (1984). An Introduction to Multivariate Statistical Analysis. Second Edition. John Wiley & Sons.
- Bai, Z. D. and Saranadasa, H. (1996). Effect of high dimension comparison of significance tests for a high dimensional two sample problem. *Statistica Sinica*. **6**, 311-329.
- Bai, Z. D. (1999). Methodologies in spectral analysis of large dimensional random matrices, A review. *Statistica Sinica*. **9**, 611-677.
- Bai, Z. D. and Silverstein, J. W. (2004). CLT for linear spectral statistics of large-dimensional sample covariance matrices. *Ann. Probab.* **32**, 553-605.
- Bai, Z. D. and Silverstein, J. W. (2006). Spectral analysis of large-dimensional random matrices, 1st ed. *Science Press, Beijing, China*.
- Dempster, A. P. (1958). A high dimensional two sample significance test. *Ann. Math. Statist.* **29**, 995-1010.
- Jonsson, D. (1982). Some limit theorems for the eigenvalues of a sample covariance matrix. *J. Multivariate Anal.* **12**, 1-38.
- Ledoit, O. and Wolf, M. (2002) Some hypothesis tests for the covariance matrix when the dimension is large compared to the sample size. *Ann. Statist.* **30**, 1081-1102.
- Pastur, L. and Lytova, A. (2008) Central limit theorem for linear eigenvalue statistics of random matrices with independent entries. *Preprint*, [arXiv:0809.4698v1\[math.PR\]](https://arxiv.org/abs/0809.4698)
- Pillai, K. C. S. (1967). Percentage points of the largest root of the multivariate beta matrix. *Biometrika*. **54**, 189-194.
- Pillai, K. C. S. and Flury, B. N. (1984). Percentage Points of the Largest Characteristic Root of the Multivariate Beta Matrix. *Communications in Statistics, Part A*. **13**, 2199-2237.
- Schott, James R. (2007). A test for the equality of covariance matrices when the dimension is large relative to the sample size. *Comput. Statist. Data Anal.* **51**, 6535-6542.
- Silverstein, J. W. (1985). The limiting eigenvalue distribution of a multivariate F-matrix. *SIAM J. Math. Anal.* **16**, 641-646.
- Srivastava, Muni S. (2005) Some tests concerning the covariance matrix in high dimensional data. *J. Japan Statist. Soc.* **35**, 251-272.

- Yin, Y. Q., Bai, Z. D. and Krishnaiah, P. R. (1983). Limiting behavior of the eigenvalues of a multivariate F-matrix. *J. Multivariate Anal.* **13**, 508-516.
- Zheng, S. (2008). Central Limit Theorem for Linear Spectral Statistics of Large Dimensional F Matrix. *Preprint, Northern-Est Normal University*

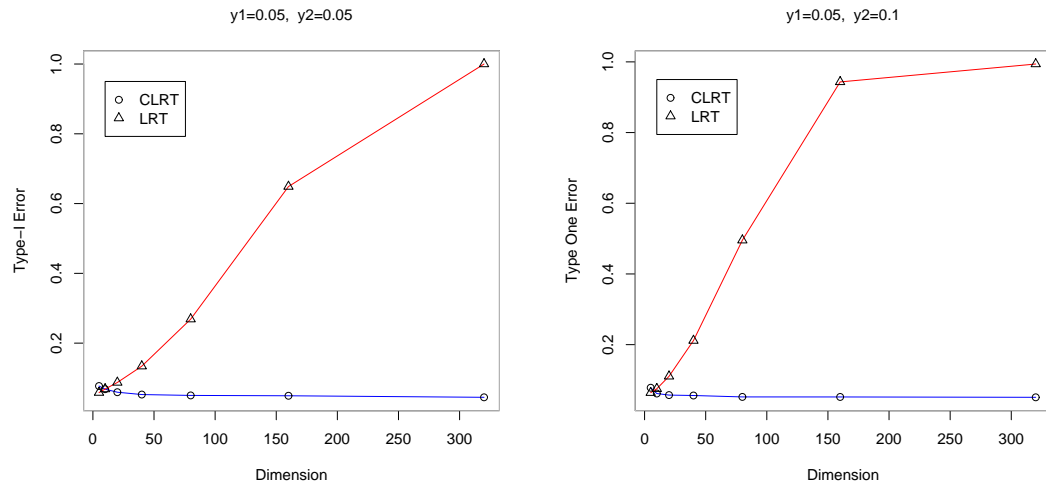


FIGURE 2. Sizes of the traditional LRT and the corrected LRT based on 10,000 independent replications using real Gaussian variables. Left: $y_1 = y_2 = 0.05$. Right: $y_1 = 0.05, y_2 = 0.1$.